



TITLE:

Bioinformaticsとソフトウェア (符号と暗号の代数的数理)

AUTHOR(S):

田辺, 隆人

CITATION:

田辺, 隆人. Bioinformaticsとソフトウェア (符号と暗号の代数的数理).
数理解析研究所講究録 2005, 1420: 28-40

ISSUE DATE:

2005-04

URL:

<http://hdl.handle.net/2433/47179>

RIGHT:

Bioinformatics とソフトウェア

(株)数理システム 田辺隆人(Takahito Tanabe)
Mathematical Systems Inc.

Cypripedium



Lady's slipper orchid

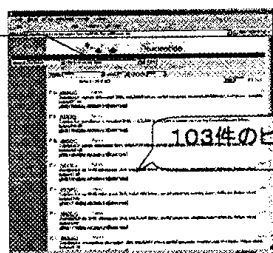
主要な配列データベース

- GenBank (<http://www.ncbi.nlm.nih.gov/>)
- EMBL (<http://www.ebi.ac.uk/embl/>)
- DDBJ (<http://www.ddbj.nig.ac.jp/>)
- PIR (<http://www-nbrf.georgetown.edu/>)
- EXPASY (<http://www.expasy.ch/>)
- SRS (<http://srs6.ebi.ac.uk/>)

データベース検索

- NCBI, ヌクレオチドデータベース

"Cypripedium
Japan"
で検索



データベース検索

- 個々のアイテムの概要

country:Canada,Prince Edward Island, isolate:AC-02
gi|13517201|dbj|AB056315.1|[13517201]

■ 4. AB056314 Reports
Cypripedium sp. U-03 chloroplast DNA, trnL(UAA) intron, partial sequence,
country:Japan,Hokkaido, Rebun Island, isolate:U-03
gi|13517200|dbj|AB056314.1|[13517200]

■ 5. AB056313 Reports
Cypripedium sp. U-02 chloroplast DNA, trnL(UAA) intron, partial sequence,

データベース検索

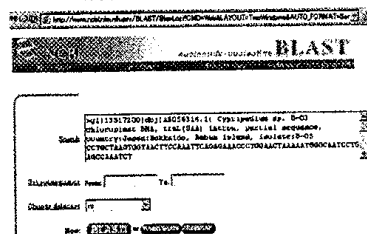
- 配列 (FASTA形式)

■ 4. AB056314 Reports Cypripedium sp. U.
gi|13517200]

```
>gi|13517200|dbj|AB056314.1| Cypripedium sp. U-03 chloroplast DNA, trnL(UAA) intron, partial
CTGCTAAGTGGTAACCTCCAAATTCAGAGAAACCCCTGGAACTAAAAATGGGCAATCCCGAGCCAAATGT
TTGTTTTTATAAAAATGCAAAATAGATAAAAGGGAAGGTCCAGAGACTCAATGGAACTGTCTTA
ACGAATGAATTTGACTAGTTAAATGGAAGATTATTCTGCAATCCATTGCAATGGAATTTGAAAGGAA
TAGAATTGAATTTGACTAGTTAAATGGAAGATTATTCTGCAATCCATTGCAATGGAATTTGAAAGGAA
AAAGTTAATGGCAGGAGATTAAGAGAGAGCTCCGATTTACATGTCATATCCAGCAATGAAA
```

アラインメント

- BLAST (類似の配列の検索)



漸化式

Needleman-Wunsch アルゴリズム

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) & \text{一致するときのスコア} \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases}$$

$$F(i, 0) = -id$$

$$F(0, j) = -jd$$

開始点をずらしたときのペナルティ

どちらかを飛ばしたときのペナルティ

インパクト

・ 進化の経路の解析

シロイヌナズナ

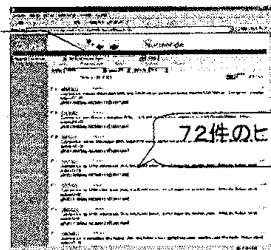
Arabidopsis thaliana

シロイヌナズナは、世代結実が早くゲノムがコンパクトなため遺伝率の高いモデル植物として広く使われ、多数の突然変異、クマクロソウ、形質転換系統等の資源が利用可能です。セントロメアやテロメアやRNAの位置配列情報を集めて、ゲノムの全塩基配列がほぼ決まっています。ゲノムデータベースの作成している基盤生物の中でも最もゲノム、全塩基配列が済んだのはシロイヌナズナです(新巻やGakopoulosではゲノムRNAメタデータが豊富です)。様々な植物生産現象の研究に使われていますが、質塩基配列にまで到達するほどゲノムメタデータの研究にも進んだ材料です。

アラインメント解析で69%の遺伝子機能が判明
(実験的には9%)

データベース検索

・ NCBI, 構造データベース

Chloroplast
で検索

72件のヒット

データベース検索

・ 個々の結果

1. 1034
Crystal Structure Of Mutant S188a Of Photosynthetic Glycerate-3-Phosphate Dehydrogenase A4 Isoform, Complexed With NADp
[mmdid28548]

Description: Crystal Structure Of Mutant S188a Of Photosynthetic Glycerate-3-Phosphate Dehydrogenase A4 Isoform, Complexed With NADp.

Deposition: F. Sparta, S. Ferman, G. F. F. A. R. P. Sabatini, P. Pupillo & P. Tosti, 27-Nov-03

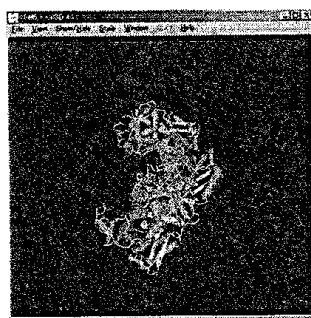
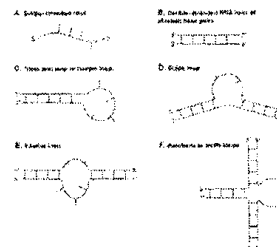
Taxonomy: [Succisa pratensis](#)

Reference: [PubMed](#) [MMDB](#) [PDB](#) [1RLH5](#)

ほうれん草

データベース検索

・ 三次構造

二次構造特定
RNAfolding

二次構造特定 入力データ

- RNA塩基配列
- 二塩基間結合のエネルギー
- 各種ループ構造のエネルギー

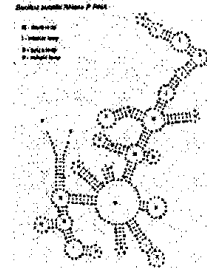
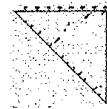
Table 3.1. Free energy of formation of RNA secondary structure at 37°C. The minimum free energy of formation of a secondary structure is given in kcal/mol.

A. Free energy of formation of base pairs		B. Free energy of formation of loops	
Base pair	Free energy (kcal/mol)	Loop type	Free energy (kcal/mol)
AA	-1.0	1	-1.0
AG	-1.0	2	-1.0
AA	-1.0	3	-1.0
AG	-1.0	4	-1.0
AA	-1.0	5	-1.0
AG	-1.0	6	-1.0
AA	-1.0	7	-1.0
AG	-1.0	8	-1.0
AA	-1.0	9	-1.0
AG	-1.0	10	-1.0

B. Destabilizing energies for loops		C. Free energy of formation of loops	
Number of bases	Free energy (kcal/mol)	Loop type	Free energy (kcal/mol)
1	-1.0	1	-1.0
2	-1.0	2	-1.0
3	-1.0	3	-1.0
4	-1.0	4	-1.0
5	-1.0	5	-1.0
6	-1.0	6	-1.0
7	-1.0	7	-1.0
8	-1.0	8	-1.0
9	-1.0	9	-1.0
10	-1.0	10	-1.0

二次構造特定 アルゴリズム

- エネルギー最小化



二次構造特定 アルゴリズム

$$W(i, j) = \min \{ W(i+1, j), W(i, j-1), V(i, j), \min_{1 \leq k < j} \{ W(i, k) + W(k+1, j) \} \} \quad (1)$$

$$V(i, j) = \min \{ V(i, j-1), V(i+1, j), V(i, j-1) + V(i+1, j-1), VM(i, j), VM(i, j-1) \} \quad (2)$$

where

$$VM(i, j) = \min_{\substack{1 \leq i' < j' < j \\ i - i' + j - j' > 2}} \{ W(i, i') + V(i', j') \} \quad (3)$$

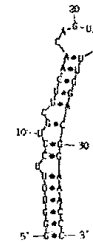
and

$$VM(i, j) = \min_{1 \leq k < j-1} \{ W(i+1, k) + W(k+1, j-1) \} \quad (4)$$

The minimum folding energy, E_{\min} , is given by $W(1, n)$.

エネルギー最小化
2-loopまで考慮

二次構造の特定 mfoldサーバー



"GGGUUUUCCUGCUCAACAG
UGCUUGGACGGAAACCC"に
対する mfoldサーバーの応答

<http://www.bielefeld.rpt.de/applications/mfold/si/rna/terml.cgi>

df = -10.1 (initially -10.1) forritin

インパクト RNA機能調整の解明

- 血中フェリチン
 - IRE (28ゲノム)にIRE結合蛋白が結合
 - mRNA安定化
 - 蛋白質を大量に合成

Examples of IREs



特徴的な構造

三次構造特定 入出力

```

1 mdsdskps ydfidliag ssaglaase aakfdkdvw idvtolpiz twagigtor
81 nvgcipkka hqaaliga: kdsnygki edtvkhdek atesvorhig sineprvel
121 rakcevyena ygfifshki: mtrckgkai vssarflia tserprylai nedecyias
181 ddfisloyap gktlvagay valacagla gldvltvaw railirgdc dmekigahn
241 ehgikfira fvkisalei agtagikrti skotnsat: edfotvilia grdootli
301 gietvyrkin aktkiercd aectvoviy aliadiakli eltpvalag zliaerlyag
361 stvkcdydv attvitpoy accalsseks vektasenie vyaffwale atvsardnek
421 cyakvienik dnervgthv igrnagvliq gfaalukgi tkqkdostig ihpocseift
481 tivskrags dliagsgoc
  
```



三次構造特定

- コンペティション



Fifth Community Wide Experiment on the Critical Assessment of Techniques for Protein Structure Prediction

Asilomar Conference Center
December 2002

三次構造特定 アルゴリズム

- 断片(3-9残基)の検索
- MonteCarlo法
- Superfamily に特徴的な構造情報利用
- ポテンシャル関数(疎水性・free-energy. .)
- 相同解析の結果を利用
- Knowledge base 検索

帰納的な方法が不可欠

機械学習 ソフトウェア

- ソフトウェア
HMMpro(隠れマルコフモデル解析ソフト)

開発元: Net-ID Inc.

- 応用
 - 遺伝子発見
 - 構造情報の取得
 - 相同解析

隠れマルコフモデル

- 記号列

- DNA/RNAアルファベット(4種)

A = Adenine C = Cytosine
G = Guanine T = Thymine/U = Uracil

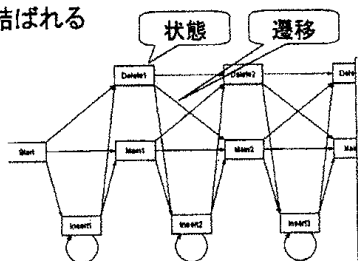
- アミノ酸(20種)

A = Alanine (Ala)	R = Arginine (Arg)
D = Aspartic Acid (Asp)	C = Cysteine (Cys)
Q = Glutamine (Gln)	E = Glutamic Acid (Glu)
N = Asparagine (Asn)	G = Glycine (Gly)
H = Histidine (His)	I = Isoleucine (Ile)
L = Leucine (Leu)	K = Lysine (Lys)
M = Methionine (Met)	F = Phenylalanine (Phe)
P = Proline (Pro)	S = Serine (Ser)
T = Threonine (Thr)	W = Tryptophan (Trp)
Y = Tyrosine (Tyr)	V = Valine (Val)

状態遷移確率 + 出力確率

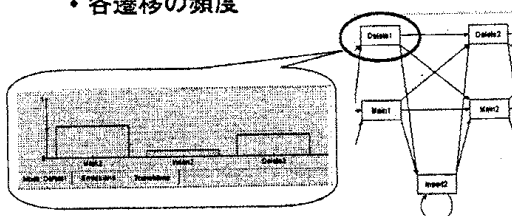
状態列

- 遷移で結ばれる



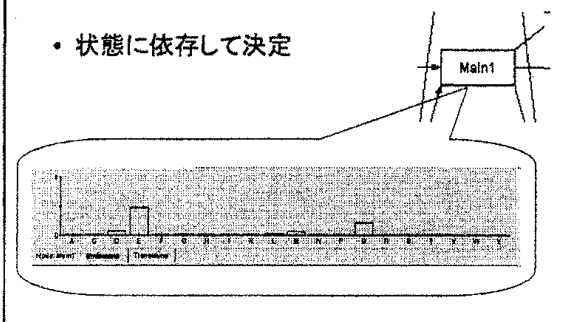
状態遷移確率

- 各遷移の頻度



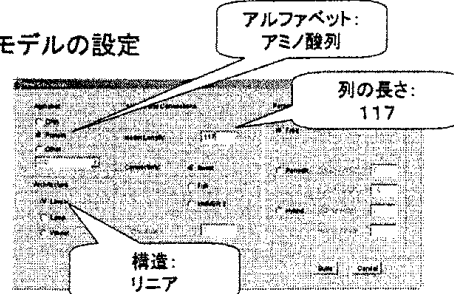
出力確率

- 状態に依存して決定



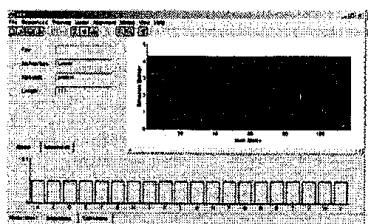
マルコフモデルの学習

- モデルの設定



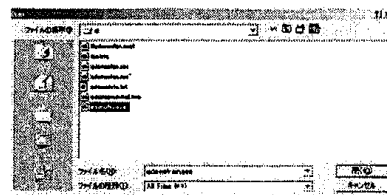
学習前

- 出力分布



学習元ファイルのインポート

- ファイル選択



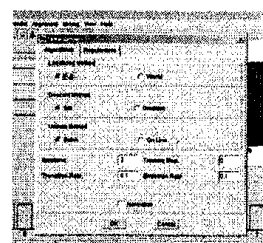
学習の実施

- 教師データ

ObjectData:
DATA: string
MKLPVRLVLMFWIPASSSDVVMTOITPLSLPVSLQDQASISCRSSOSLVHSHGNTYLNWYLO
KAGOSPKLLVYKVSNIPOVDFRPSGSGSOTDFTLKISKVBAEDLGIFCSQTHVPTFGGOT
KLEIKR
LOOPGAEVLKPGASVILCKASGYTFITNYWVWVKORPGRGLEWIGRIDPNSGGTKYNEKF
KNAATLTINKPNTAYMQLSSLTDDSAVYYCARGYDYSYAMDYWGQOTSVTVSS
ESGGGLVQPGGSMKLSVCASGPTTSNYWMNVYROSPEKLEWVAERLKSGYATHYAESV
KGRFTSRDDSSVYLOMNNLKAEDTGIYYCTRPVDPDYWGQOTTLTVSS
MEPOLSWIELVAILKGOVCEVRLVESGGDLVEPGSLXVSCVSGPISKAWMNVYRQAPG
KOLQWVGOKNHYDQGTIDYAAVVGGRPERDDSKSTVYLQMNELKIEDTAVYYCYQNT
GTVDYWGGQTLTVTVSS
ISCKASGYTFITNYGMNVVVKQAPGKGLKWMGWINTYTGEPYADDFKGRFAFSLTSASTA
YLONNLKNEDTATPCARGSYDYIYAMDYWGQOTSVTVSS
LVOLQSGQPLVKPGTSMKISCKTSYSGFTGYTMSWVROSHGKSLWIGLIPSGGTNYNO
KPKDKASLTVDKSSSTAYMELLSLTSDDSAVYYCARPSYYSGRNYYAMDYWGQOTSVTVS
SAK

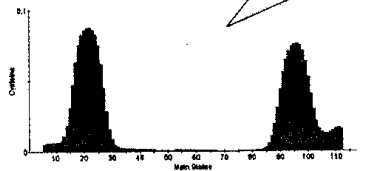
学習の実施

- アルゴリズム・パラメータ設定



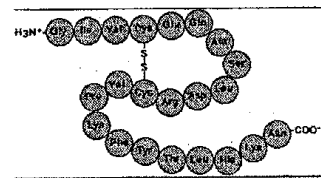
学習結果

- 学習結果 (Cysteine)



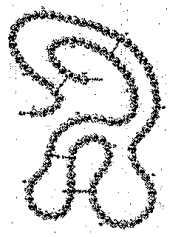
Disulfide bridge?

- 概念図



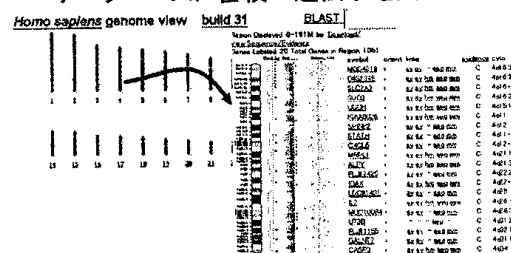
Disulfide bridge

- Lysozyme (4箇所出現)

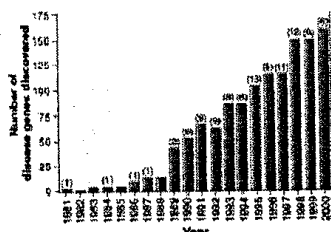


インパクト 遺伝子発見

- データベースに蓄積→遺伝子地図



インパクト 医療方面への応用



計算機科学への要請

- インタフェース
- 最適化
- 全文検索
- 圧縮

— Ugly programming instead of
deep biological insights

全文検索 索引作成

- サフィックスとインデックスポイント

Text	a	b	r	a	c	a	d	a	b	r	a
Index	0	1	2	3	4	5	6	7	8	9	10
Suffix	Index										
a	b	r	a	c	a	d	a	b	r	a	
b	r	a	c	a	d	a	b	r	a		0
r	a	c	a	d	a	b	r	a			1
a	c	a	d	a	b	r	a				2
c	a	d	a	b	r	a					3
d	a	b	r	a							4
a	b	r	a								5
b	r	a									6
r	a										7
a											8
											9
											10

全文検索 索引作成

- ソート

Sorted Suffix	Index
a	10
a b r a	7
a b r a c a d a b r a	0
a c a d a b r a	3
a d a b r a	5
b r a	8
b r a c a d a b r a	1
c a d a b r a	4
d a b r a	6
r a	9
r a c a d a b r a	2

全文検索 索引の利用

- 二分探索

Sorted Suffix	Index
a	10
a b r a	7
a b r a c a d a b r a	0
a c a d a b r a	3
a d a b r a	5
b r a	8
b r a c a d a b r a	1
c a d a b r a	4
d a b r a	6
r a	9
r a c a d a b r a	2

全文検索 索引作成アルゴリズム

- ブロックソートに帰着

入力文字列 + \$

Suffix array
- \$abcabc
0 abcabc\$
1 bcabc\$a
2 cabc\$bc
3 abc\$abc
4 bc\$abca
5 c\$abcab

ソート

- \$abcabc
3 abc\$abc
0 abcabc\$
4 bc\$abca
1 bcabc\$a
5 c\$abcab
2 cabc\$bc

\$: いずれの文字よりも小

ブロックソート

入力文字列

0 AbcABC	3 ABCAbc
1 bcABCA	0 AbcABC
2 cABCAb	4 BCAbcA
3 ABCAbc	5 CAbcAB
4 BCAbcA	1 bcABCA
5 CAbcAB	2 cABCAb

ソート

入力が特殊なソート

ブロックソートによる変換 Burrows-Wheeler 変換

- 着眼

入力文字列

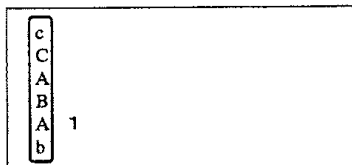
0 AbcABC	3 ABCAbc
1 bcABCA	0 AbcABC
2 cABCAb	4 BCAbcA
3 ABCAbc	5 CAbcAB
4 BCAbcA	1 bcABCA
5 CAbcAB	2 cABCAb

ソート

最終列と1番行の場所で全情報再現可能

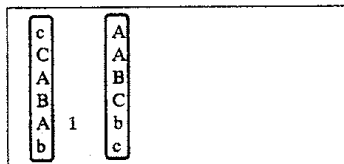
Burrows-Wheeler 変換 逆変換

- 入力



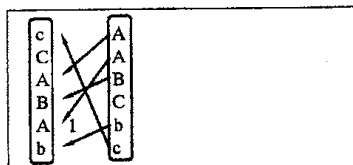
Burrows-Wheeler 変換 逆変換

- ソートによって最初の行を知る



Burrows-Wheeler 変換 逆変換

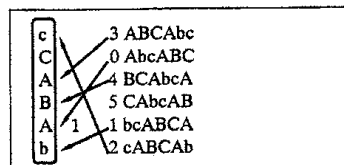
- 行番号特定



(i+1) 番行の最終文字
= i 番行の先頭文字

Burrows-Wheeler 変換 逆変換

- Suffix array 及び 元の文字列判明



suffix array ・元の文字列判明

Burrows-Wheeler 変換 メリット

```
t: hat acts like this:<13><10><1
t: hat buffer to the constructor
t: hat corrupted the heap, or wo
W: hat goes up must come down<13
t: hat happens, it isn't likely
w: hat if you want to dynamical
t: hat indicates an error.<13><1
t: hat it removes arguments from
t: hat looks like this:<13><10><
t: hat looks something like this
t: hat looks something like this
t: hat once I detect the mangled
```

繰り返しが見えやすい→符号化効率良好

圧縮と検索の両立

ブロックソート



BWTによる高性能圧縮

+

suffix array による高速全文検索

ブロックソート アルゴリズム

- quick sort
一般アルゴリズムなので向かない
- ternary partitioning[Bentley, Sedgewick 97]
無駄な文字列比較が少ない
- doubling algorithm
多くの場合最速
- copy method
対象の性質を利用
- layer method
copy method の改良

’04 田辺・小林

copy method 原理

- 先頭二文字のみでソートし、「区間」決定

```

{ aa.....
{ ab.....
{ ab.....
{ ab.....
{ ab.....
{ ae.....
{ ax.....
{ az.....
{ az.....
{ ba.....
  
```

copy method 原理

- 先頭文字が同一の「区間」内を順にソート

```

{ aa.....
{ ab.....
{ ab.....
{ ab.....
{ ab.....
{ ae.....
{ ax.....
{ az.....
{ az.....
{ ba.....
  
```

copy method 原理

- ソートされた区間内の末尾文字を調べる

```

{ aa.....h
{ ab.....f
{ ab.....e
{ ab.....u
{ ab.....f
{ ae.....g
{ ax.....
{ az.....
{ az.....
{ ba.....
  
```

copy method 原理

- 判明箇所を決定(コピー)

```

{ aa.....h
{ ab.....f
{ ab.....e
{ ab.....u
{ ab.....f
{ ae.....g
{ ax.....
{ az.....
{ az.....
{ ba.....
  
```

一番速いのは「a」区間の最初

一番速いのは「a」区間の二番目

copy method 原理

- 判明箇所を決定(コピー)

```

{ ex.....
{ ex.....
{ fab.....
{ fab.....
{ fa.....
{ fa.....
{ fa.....
{ fb.....
{ fb.....
{ fb.....
  
```

コピー先の区間では上から決定される

